

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-189502

(43)Date of publication of application : 05.07.2002

(51)Int.Cl.

G05B 13/02

G05B 13/04

G06N 3/00

(21)Application number : 2000-386265

(71)Applicant : JAPAN SCIENCE & TECHNOLOGY  
CORP  
ADVANCED TELECOMMUNICATION  
RESEARCH INSTITUTE  
INTERNATIONAL

(22)Date of filing : 20.12.2000

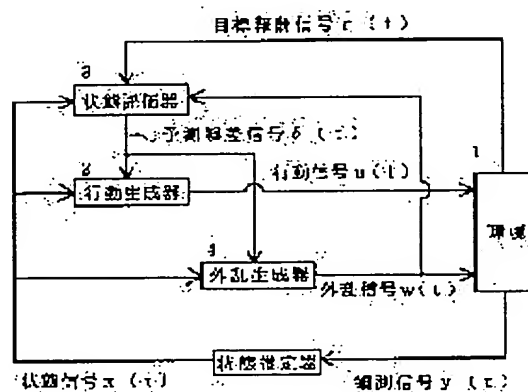
(72)Inventor : MORIMOTO ATSUSHI  
DOTANI KENJI

## (54) ROBUST REINFORCE LEARNING SYSTEM

## (57)Abstract:

PROBLEM TO BE SOLVED: To provide a method for learning robust control which is relatively resistant to the fluctuation of an environment capable of complying even with the case that an environment model is unknown and to provide a robust controller.

SOLUTION: This system is provided with an action generator (2) having a learning function and also outputting an action signal  $u(t)$  to an environment (1), a disturbance generator (4) having a learning function and also outputting a disturbance signal  $w(t)$  to an environment, and a state evaluator (3) generating an evaluation signal  $q(t)$  being a reward signal obtained by adding a reward corresponding to the level of achievement of a target with a reward corresponding to resistance to disturbance from the disturbance generator, and predicting the expected value of the sum of the evaluation signals to be obtained from the present state  $x(t)$  to the future, and generating the prediction error signal. Then, the action generator learns to maximize the expected value of the sum of the evaluation signals to be obtained from the present state to the future, while the disturbance generator learns to minimize the expected value of the sum of the valuation signals.



## LEGAL STATUS

[Date of request for examination]

20.12.2000

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3465236

[Date of registration] 29.08.2003

[Number of appeal against examiner's decision  
of rejection]

[Date of requesting appeal against examiner's  
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号  
特開2002-189502  
(P2002-189502A)

(43) 公開日 平成14年7月5日 (2002.7.5)

(51) Int.Cl. <sup>7</sup>	識別記号	F I	テ-マ-ト* (参考)
G 0 5 B 13/02		G 0 5 B 13/02	T 5 H 0 0 4
	13/04		J
G 0 6 N 3/00	5 5 0	G 0 6 N 3/00	5 5 0 E

審査請求 有 請求項の数 9 O L (全 13 頁)

(21) 出願番号 特願2000-386265 (P2000-386265)

(22) 出願日 平成12年12月20日 (2000. 12. 20)

特許法第30条第1項適用申請有り

(71) 出願人 396020800

科学技術振興事業団

埼玉県川口市本町4丁目1番8号

(71) 出願人 393031586

株式会社国際電気通信基礎技術研究所

京都府相楽郡精華町光台二丁目2番地2

(72) 発明者 森本 淳

奈良県生駒市高山町8916-5 学生宿舎2-102

(74) 代理人 100099265

弁理士 長瀬 成城

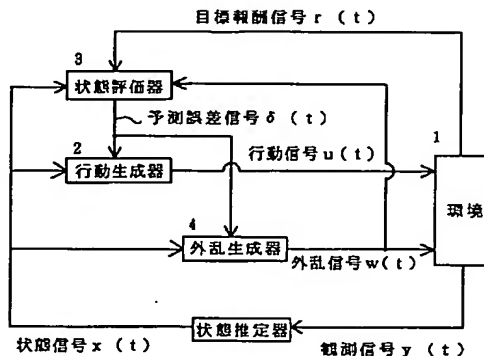
最終頁に続く

(54) 【発明の名称】 ロバスト強化学習方式

(57) 【要約】

【課題】 環境モデルが未知の場合にも対応することができるとともに、環境の変動に比較的強いロバスト制御の学習方法およびロバスト制御器を提供する。

【解決手段】 環境 (1) に行動信号  $u(t)$  を出力するとともに学習機能を具備する行動生成器 (2) と、環境に外乱信号  $w(t)$  を出力するとともに学習機能を具備する外乱生成器 (4) と、目標の達成度に応じた報酬に、前記外乱生成器からの外乱に耐えうることに応じた報酬を加味した報酬信号である評価信号  $q(t)$  を生成し、現在の状態  $x(t)$  から将来に向けて得られる評価信号の和の期待値を予測し、その予測誤差信号を生成する状態評価器 (3) とを備え、現在の状態から将来に向けて得られる評価信号の和の期待値を最大化するべく行動生成器は学習し、一方、外乱生成器は前記評価信号の和の期待値を最小化すべく学習する。



## 【特許請求の範囲】

【請求項1】制御対象あるいは環境に行動信号を出力するとともに学習機能を具備する行動生成器、および制御対象あるいは環境に外乱信号を出力するとともに学習機能を具備する外乱生成器を備え、

目標の達成度に応じた報酬に、前記外乱生成器からの外乱に耐えることに応じた報酬を加味した報酬信号である評価信号を生成し、現在の状態から将来に向けて得られる評価信号の荷重和の期待値を最大化（または最小化）するべく行動生成器は学習し、一方、外乱生成器は前記評価信号の和の期待値を最小化（または最大化）するべく学習することを特徴とするロバスト強化学習方式。

【請求項2】前記学習方式において、現在の状態から将来に向けて得られる評価信号の和の期待値を予測する状態評価器を備え、その予測誤差信号を、状態評価器、行動生成器、および外乱生成器の少なくとも1個の学習に用いることを特徴とするロバスト強化学習方式。

【請求項3】前記状態評価器、行動生成器および外乱生成器の少なくとも一個は、関数近似手段として、入出力関係を示す参照テーブルを具備していることを特徴とする請求項1または2に記載のロバスト強化学習方式。

【請求項4】前記状態評価器、行動生成器および外乱生成器の少なくとも一個は、関数近似手段として、線形モデルまたは多項式モデルを具備していることを特徴とする請求項1または2に記載のロバスト強化学習方式。

【請求項5】前記状態評価器、行動生成器および外乱生成器の少なくとも一個は、関数近似手段として、多層神経回路網を具備していることを特徴とする請求項1または2に記載のロバスト強化学習方式。

【請求項6】請求項1または請求項2の方式により、予め学習された前記状態評価器と行動生成器または行動生成器のみを用いた制御方式。

【請求項7】請求項1または請求項2の方式を計算機シミュレーションによって実現される環境モデルに適用し、それによって学習された前記状態評価器と行動生成器または行動生成器のみを実環境に適用することを特徴とする請求項6に記載の制御方式。

【請求項8】前記状態評価器または行動生成器の少なくとも一方は、関数近似手段として、入出力関係を示す参照テーブルを具備していることを特徴とする請求項6または7に記載のロバスト制御方式を用いたロバスト制御器。

【請求項9】前記状態評価器または行動生成器の少なくとも一方は、関数近似手段として、線形モデル、多項式

モデルまたは多層神経回路網を具備していることを特徴とする請求項6または7に記載のロバスト制御方式を用いたロバスト制御器。

## 【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、強化学習において環境の不確かさに対してロバストな行動則を学習する方法に関するものである。さらに、詳しくは、ロボット、自動車、航空機などの物理的システムの制御、また、広く人間に代わって情報検索、ユーザー応答、資源割当、市場取引などを行うコンピュータプログラムなどにおいて、環境の様々な外乱や、想定していたモデルと現実との環境のズレに対応することができる強化学習方式および強化学習された制御器である。なお、この明細書においては、特許明細書では使用不可能な文字に対応するために、下記のようにして対処している。

【外1】

$H_{\infty}$  は  $H^{\infty}$  で表している。

$\int_{\cdot}^{\infty}$  は  $\int t^{\infty}$  で表している。

【0002】

【従来の技術】従来の強化学習方式について説明する。図2は従来の学習方式に用いる回路のブロック図である。図3は制御対象と制御器とのフィードバックのブロック図であり、 $H_{\infty}$ 制御理論を説明するための図である。環境1は制御対象をはじめ、行動の対象となるシステム一般を表し、状態信号  $x(t)$  を行動生成器2および状態評価器3に出力している。行動生成器2は状態信号  $x(t)$  が入力されると行動信号  $u(t)$  を環境1に出力している。また、環境1から報酬信号  $r(t)$  が状態評価器3に入力される。状態評価器3は、目標の達成度に応じた報酬  $r(t)$  が入力されると、現在の状態  $x(t)$  から将来に向けて得られる、報酬  $r(t)$  の荷重和の期待値すなわち、評価関数  $V(x(t))$  を予測し、その予測値を用いて予測誤差信号  $\delta(t)$  を生成し、行動生成器2に出力する。行動生成器2は、状態評価器3から予測誤差信号  $\delta(t)$  が入力されると、現在の状態  $x(t)$  から将来に向けて得られる、報酬  $r(t)$  の荷重和の期待値すなわち、評価関数  $V(x(t))$  が最大となる様に学習し、その入出力の関係を変更する。ただし、前記評価関数  $V(x(t))$  は、連続の場合、

【数1】

$$V(x(t)) = E \left[ \int_t^{\infty} e^{-\frac{s-t}{\tau}} r(x(s), u(s)) ds \right]$$

(式1)

(τは評価の時定数)、離散系の場合

$$V(x_T) = E \left[ \sum_{s=t}^{\infty} \alpha^{s-t} r_s \right]$$

(αは評価の減衰率)

と定義される。

この学習方式では、ある環境のもとで最適な行動が学習されるが、異なる環境では動作は保証されていない。また、異なる環境に適応するためには再学習を行う必要があり、その再学習の時間が新たに必要となる。

【0003】次に、従来のH<sup>∞</sup>制御について説明する。図3において、制御対象Gから観測信号y(t)が制御器Kに入力される。制御器Kは観測信号y(t)が入力されると、制御信号である行動信号u(t)を制御対象Gに出力する。また、制御対象Gには外乱信号w(t)および行動信号u(t)が入力され、これらの信号が入力されると、評価用信号z(t)および観測信号y(t)を出力する。なお、外乱の影響を評価するための評価用信号z(t)と、制御対象Gを観測して制御器Kに入力するフィードバック信号である観測信号y(t)とは、同じにする\*

\*ことも可能であるが、異ならしめることも可能である。

そして、ロバスト制御の代表的な定式化であるH<sup>∞</sup>制御問題の要請は、図3に図示するフィードバック系で未知外乱やモデル誤差に起因する外乱信号w(t)による評価用信号z(t)への影響を少なく抑えつつ、出力を安定化する。すなわち、評価用信号z(t) = 0に近づけることである。具体的には、H<sup>∞</sup>ノルムによりシステムの外乱に対する感度を測り、ロバスト性の基準値γ以下となるような制御器Kの設計を行う。ノルムとは、ある種の大きさの指標であり、外乱信号w(t)から評価用信号z(t)への伝達関数行列をT<sub>zw</sub>としたとき、そのH<sup>∞</sup>ノルム、||T<sub>zw</sub>||<sub>∞</sub>は次の(式2)のように定義される。

【数2】

$$\|T_{zw}\|_{\infty} = \sup_w \frac{\|z\|_2}{\|w\|_2} < \gamma \quad (\text{式2})$$

ただし、sup<sub>w</sub>は外乱信号w(t)に関する上限を表し、30※それぞれ評価用信号z(t)および外乱信号w(t)のL<sub>2</sub>ノルムであり、次の(式3)および(式4)で定義される。外乱信号w(t)を変化させたときに、||z||<sub>2</sub> / ||w||<sub>2</sub>がsup<sub>w</sub>(||z||<sub>2</sub> / ||w||<sub>2</sub>)より大きくなならないことを示している。また、||z||<sub>2</sub>、および||w||<sub>2</sub>はそれぞれ

【数3】

$$\|z\|_2 = \left( \int_0^{\infty} z^T(t) z(t) dt \right)^{\frac{1}{2}} \quad (\text{式3})$$

【数4】

$$\|w\|_2 = \left( \int_0^{\infty} w^T(t) w(t) dt \right)^{\frac{1}{2}} \quad (\text{式4})$$

【0004】ここで、評価関数Vを次の(式5)の様に★【数5】  
★  
定義する。

$$V = \int_0^{\infty} (-z^T(t) z(t) + \gamma^2 w^T(t) w(t)) dt$$

(式5)

これを行動信号u(t)に関しては最大化し、外乱信号w(t)に対しては最小化する問題を考える。その結果、V ≥ 0を満たす解を得られれば、(式2)の条件のもと

で、評価用信号z(t)の安定化が実現できる。

【0005】強化学習の課題のうち、予め与えられた目標点あるいは目標軌道への近さを報酬信号とするもの

は、学習制御の課題と考えることができる。前記ロバスト制御の代表的な方法である $H^\infty$ 制御は、システムの外乱による影響の受けやすさを $H^\infty$ ノルムで評価し、フィードバック系の $H^\infty$ ノルムを一定以下に抑える制御器を設計することにより、外乱やモデル誤差に対するロバスト性を保証するものである。しかし、その制御器の解析的な構成手法は線形システムに対し示されており、非線形システムに対してはある限定されたシステムを除いては、一般に解析的に制御器を構築する方式はない。非線形システムにおいて、多層神経回路網を用いて未知外乱を考慮した状態価値関数を近似しロバスト制御を実現する手法が提案されているが、これらは、制御器の適応可能範囲が線形近似可能な領域付近に限られていたり、学習がオフラインのバッチ学習に限られている。また、これらの学習には環境のモデルを必要としている。さらに、状態を離散化し動的計画法を用いる方式も提案されているが、制御器の構築には、状態を離散化する過程と、オフラインの計算過程とを必要とし、かつ、環境モデルを必要とする。また、ここまで挙げたロバスト制御器はレギュレータ（目標点を原点とし、その原点に制御対象の状態を持って行く制御）としてのみ機能する。

【0006】ところで、従来のミニマックス戦略を用いる強化学習方式では、オセロやバックギャモンの学習を行っている。また、戦闘機やミサイルの逃亡追従問題に適用した例もある。一方、本発明では、仮想的な敵を想定して学習を行った後、実際に制御器を適用する際には、その外乱生成器を除いて制御器を構成するという点で異なる。これは、外乱を仮想的な敵とみなすことによりロバストな制御器を強化学習により構築するという新しい発想に依存している。また、後述の（式8-2）で定義される評価関数の予測誤差に重み付けをして学習することで、ロバストな行動則を強化学習によって獲得する手法とは異なり、本発明は $H^\infty$ 制御との関連や目的関数が明確である。よって、最悪の外乱を求めているという点において、より確実にロバストな制御器を獲得することができる。

【0007】そして、後述のように、本発明の実施の形態においては、非線形、オンライン、モデル非依存性の動的最適化手法である強化学習方式に、ロバスト制御における最悪外乱に対する感度の最小化の原理を導入することにより、ロバスト強化学習方式を実現し、その制御課題への適用により、非線形、オンライン、モデル非依存性のロバスト学習制御を実現する。また、環境の近似モデルが既知の場合、あるいは未知の場合でも、環境のモデルを学習することにより、モデル依存の強化学習方式を拡張し適用することにより、効率良く非線形ロバスト制御器をオンラインで構築することができる。

【0008】この様なことができる利点は、環境が非定常であった場合に、一度学習によって構築した非線形ロバスト制御器や環境モデルを用いれば、ある程度の環境

の変動なら再学習することなく対応することができることである。また、従来のロバスト制御は出力をある目標値に収束させるレギュレータ問題に対してのみ定式化されていたが、本発明のロバスト強化学習方式は、任意の評価基準に対して適用可能である。

【0009】

【発明が解決しようとする課題】この様に、従来の強化学習では、環境の変動が生じた場合には、再学習する必要があり、従来のロバスト制御の方式は環境のモデルが未知の場合には対応することができなかった。さらに、非線形系に適用可能な制御則の具体的方法は与えていない。

【0010】そこで本発明は、環境モデルが未知の場合にも対応することができるとともに、非線形系に適用可能な環境の変動に強いロバスト強化学習方式およびロバスト制御器を提供することを目的とする。

【0011】

【課題を解決するための手段】このため、本発明が採用した課題解決手段は、制御対象あるいは環境に行動信号を出力するとともに学習機能を具備する行動生成器、および制御対象あるいは環境に外乱信号を出力するとともに学習機能を具備する外乱生成器を備え、目標の達成度に応じた報酬に、前記外乱生成器からの外乱に耐えうることに応じた報酬を加味した報酬信号である評価信号を生成し、現在の状態から将来に向けて得られる評価信号の荷重和の期待値を最大化（または最小化）するべく行動生成器は学習し、一方、外乱生成器は前記評価信号の荷重和の期待値を最小化（または最大化）すべく学習することを特徴とする強化学習方式であり、前記学習方式において、現在の状態から将来に向けて得られる評価信号の和の期待値を予測する状態評価器を備え、その予測誤差信号を、状態評価器、行動生成器、および外乱生成器の少なくとも1個の学習に用いることを特徴とする強化学習方式であり、前記状態評価器、行動生成器および外乱生成器の少なくとも1個は、関数近似手段として、入出力関係を示す参照テーブルを具備していることを特徴とする強化学習方式であり、前記状態評価器、行動生成器および外乱生成器の少なくとも1個は、関数近似手段として、線形モデルまたは多項式モデルを具備していることを特徴とする強化学習方式であり、前記状態評価器、行動生成器および外乱生成器の少なくとも1個は、関数近似手段として、多層神経回路網を具備していることを特徴とする強化学習方式であり、前記方式により、予め学習された前記状態評価器と行動生成器または行動生成器のみを用いた制御方式であり、前記方式を計算機シミュレーションによって実現される環境モデルに適用し、それによって学習された前記状態評価器と行動生成器または行動生成器のみを実環境に適用することを特徴とする制御方式であり、前記状態評価器または行動生成器の少なくとも一方は、関数近似手段として、入出力関

係を示す参照テーブルを具備していることを特徴とするロバスト制御器であり、前記状態評価器または行動生成器の少なくとも一方は、関数近似手段として、線形モデル、多項式モデルまたは多層神経回路網を具備していることを特徴とするロバスト制御器である。

【0012】

【実施の形態】次に、本発明における強化学習方式およびロバスト制御器の実施の一形態を説明する。図1は本発明の実施の形態のロバスト強化学習方式に用いる回路のブロック図である。図4は第1具体例の説明図で、

(a)が概略図、(b)が1mの長さの振り子の角度変化のグラフである。図5は第1具体例の角度変化のグラフで、(a)が0.5mの長さの振り子のグラフ、

(b)が2mの長さの振り子のグラフである。図6は第2具体例の概略図である。図7は第2具体例の位置変化のグラフで、(a)が1kgの質量の搬送物を用いた場合のグラフ、(b)が3kgの質量の搬送物を用いた場合のグラフである。

【0013】本発明では、外乱や環境の変化に対してロバストな強化学習を実現するため、目標の達成度に応じた報酬 $r(t)$ に対して、外乱に耐えうることに応じた報酬 $s(t)$ を加えた新たな報酬である評価信号 $q(t)$ を次のように定義した時、

【0014】

【数6】

$$q(t) = r(x(t), u(t)) + s(w(t)) \quad (\text{式6})$$

この評価信号 $q(t)$ を報酬とした最大最小問題を強化学習方式の枠組みにおいて解く。よって、前述の $H^\infty$ 制御の問題設定は、本発明の実現例の一つとなっていることがわかる。以上を考慮した上で、次のようにロバスト\*

$$V(x(t)) = E \left[ \int_t^\infty e^{-\frac{\tau-t}{\tau}} q(s) ds \right]$$

(式7)

行動生成器2および外乱生成器4は、この期待値 $V(x(t))$ がそれぞれ、最大化、最小化される様に行動信号 $u(t)$ 、外乱信号 $w(t)$ を学習する。なお、行動生成器2、状態評価器3および外乱生成器4としては、参照テーブル、線形モデル、多項式モデル、多層神経回路網などを用いることができる。

【0017】ロバスト強化学習を行う時点においては、状態評価器3、行動生成器2および外乱生成器4は同時に作動させるが、実際に学習した行動則を制御対象または環境に用いる段階においては、状態評価器3および行動生成器2、或いは、行動生成器2のみを用いて動作させる。この行動生成器2には、観測信号 $y(t)$ として状態信号 $x(t)$ が直接得られる場合はそれを用いるが、一般にはオブザーバ、カルマンフィルタなどにより状態信

\*強化学習の学習方式に用いる回路を図1に図示するように構築する。

【0015】この図1の説明において、図2の従来の回路と同じ構成要素には同じ符号を付して、その説明は省略する。この図1においては、外乱生成器4が設けられている。そして、環境1からの観測信号 $y(t)$ が、状態推定器を介して状態信号 $x(t)$ となり、行動生成器2、状態評価器3および外乱生成器4に入力されている。この外乱生成器4は状態信号 $x(t)$ が入力されると、外乱信号 $w(t)$ を環境1および状態評価器3に出力する。この状態評価器3は、環境1からの目標報酬信号 $r(t)$ と、外乱生成器4からの外乱信号 $w(t)$ に基づいて生成した外乱報酬信号 $s(t)$ とに基づいて、現在の状態 $x(t)$ から将来に向けて得られる評価信号 $q(t)$ の荷重和の期待値を予測し、その予測値に基づいて予測誤差信号 $\delta(t)$ を生成し、行動生成器2および外乱生成器4に出力する。この様にして、状態評価器3は、外乱報酬信号 $s(t)$ に目標報酬信号 $r(t)$ を加算して評価信号 $q(t)$ を得て、予測誤差信号 $\delta(t)$ を生成し出力している。そして、行動生成器2は予測誤差信号 $\delta(t)$ が入力されると、現在の状態 $x(t)$ から将来に向けて得られる上記評価信号 $q(t)$ の荷重和の期待値を最大化するべく学習し、その入出力関係を変更する。一方、外乱生成器4は予測誤差信号 $\delta(t)$ が入力されると、現在の状態 $x(t)$ から将来に向けて得られる上記評価信号 $q(t)$ の荷重和の期待値を最小化するべく学習し、その入出力関係を変更する。

【0016】状態評価器3は、(式7)で定義される現在の状態 $x(t)$ から将来に向けて得られる評価信号 $q(t)$ の期待値 $V(x(t))$ を予測する。ただし、 $\tau$ は評価の時定数である。

【数7】

号 $x(t)$ を推定し入力として用い、また、学習時には、環境1はモデルでも、実際の環境でも可能である。そして、実際の環境の場合には、行動信号 $u(t)$ および外乱信号 $w(t)$ は、アクチュエータなどの駆動源や、低レベルの制御プログラムへの指令などの作動手段を介して環境1に入力される。一方、報酬信号 $r(t)$ や状態信号 $x(t)$ は、センサーなどの検知手段を介して環境1から出力される。

【0018】そして、状態評価器3は、評価関数 $V(x(t))$ のパラメータ $v = \{v_1, v_2, \dots, v_1, \dots\}$ を持つ近似器 $V(x(t); v)$ として実現され、その手段としては、前述の参照テーブル、線形モデル、多項式モデルおよび多層神経回路網を用いることができる。

【0019】この様にして、状態評価器3は、環境1か

ら目標報酬信号  $r(t)$  を得る手段と、外乱生成器4から外乱信号  $w(t)$  を得る手段と、目標報酬に外乱報酬を加味した評価信号  $q(t)$  を得る手段と、現在の状態から将来に向けて得られる評価信号  $q(t)$  の和の期待値を予測し、予測誤差信号  $\delta(t)$  を生成する手段とを有している。また、行動生成器2は、環境1から状態信号  $x(t)$  を得る手段と、状態評価器3から予測誤差信号  $\delta(t)$  を得る手段と、環境1に行動信号  $u(t)$  を出力する手段と、現在の状態から将来に向けて得られる評価信号  $q(t)$  の和の期待値が最大化するように学習する手段とを有している。そして、外乱生成器4は、環境1から状態信号  $x(t)$  を得る手段と、状態評価器3から予測誤差信号  $\delta(t)$  を得る手段と、環境1に外乱信号  $w(t)$  を出力する手段と、現在の状態から将来に向けて得られる評価信号  $q(t)$  の和の期待値が最小化する様に学習する手段とを有している。

【0020】以降、離散系での評価関数の学習、連続系での評価関数の学習、離散系での行動決定方法、連続系での行動決定方法の順に示す。離散系での評価関数の学習：次の様な確率分布  $P$  にしたがう動的制御対象を考える。

$P(x_{t+1} | x_t, u_t, w_t)$   
ただし、 $x_t$  は状態変数、 $u_t$  は制御入力、 $w_t$  は外乱入力を表す。このとき、求めるべき評価関数  $V_t$  は次の式で表される。

【数8】

$$V_t = \sum_{s=t}^{\infty} \alpha^{s-t} q_s \quad (\text{式8})$$

ただし、 $q_t$  はただちに得られる評価信号、 $\alpha$  ( $0 \leq \alpha \leq 1$ ) は評価の減衰率を表す。そこで、状態評価値の予測誤差  $\delta_t$  は次式のように表される。

$$\delta_t = q_t + \alpha V_{t+1} - V_t \quad (\text{式8-2})$$

【0021】よって、この予測誤差  $\delta_t$  と、次の式(式20★で、9)で表される各パラメータの寄与度の履歴  $e_{i,t}$  を用いて★

$$e_{i,t} = \sum_{n=0}^T (\alpha \lambda)^{T-n} \frac{\partial V_t}{\partial v_i} \quad (\text{式9})$$

パラメータの更新量  $\Delta v_i$  は、次式のように表される。  
 $\Delta v_i = \eta \delta_t e_{i,t}$

ただし、 $\lambda$  はパラメータの寄与度の履歴の減衰率を、 $\eta$  ☆

☆は学習率を表す。また、各パラメータの寄与度の履歴  $e_{i,t}$  は(式9)の定義より次式を用いて更新される。

【数10】

$$e_{i,t} = \alpha \lambda e_{i,t-1} + \frac{\partial V_t}{\partial v_i} \quad (\text{式10})$$

【0022】連続系での評価関数の学習：次の様な動的制御対象を考える(状態変数  $x(t)$  の時間変化  $dx/dt$  を、状態変数  $x(t)$ 、制御入力  $u(t)$ 、外乱入力  $w(t)$ 、ノイズ入力  $n(t)$  の関数として考える)。

$$dx/dt = f(x(t), u(t), w(t)) + n(t)$$

ただし、このとき、求めるべき評価関数  $V(t)$  は次式で表される。

【数11】

$$V(t) = E \left[ \int_t^{\infty} e^{-\frac{s-t}{\tau}} q(s) ds \right]$$

(式11)

ただし、 $q(t)$  はただちに得られる報酬であり、 $\tau$  は評価値の時定数である。よって、状態評価値の予測誤差  $\delta(t)$  は次式のように表される。

40\* 【0023】ここで得られる状態評価値の予測誤差  $\delta(t)$  と、次式で表される各パラメータの寄与度の履歴を用いて、

$$\delta(t) = q(t) - (1/\tau) \times V(t) + dV(t)/dt \quad (\text{数12})$$

$$e_{i,t}(t) = \int_0^t e^{-\frac{t-s}{k}} \frac{\partial V(s)}{\partial v_i} ds \quad (\text{式12})$$

ただし、 $k$  はパラメータの寄与度の履歴の時定数である。状態評価器のパラメータの更新量  $dv_i/dt$  (連続系ではパラメータ  $v_i$  の時間微分で表される) は次式のように表される。

$$dv_i/dt = \eta \delta(t) e_{i,t}(t)$$

ただし、 $\eta$  は学習率を表す。

【0024】また、各パラメータの寄与度の履歴  $e_{i,t}(t)$  の更新量  $de_{i,t}(t)/dt$  は、(式12)の定義

50



により次式を用いて更新される。

$$\frac{d e_i(t)}{d t} = -\frac{1}{k} e_i(t) + \frac{\partial V(t)}{\partial v_i} \quad (\text{式 13})$$

【0025】離散系での行動決定方法：

※て、状態  $s$  における行動  $a$  を決定する。

(モデル非依存の場合) 次式で示す確率分布  $P_r$  に従っ※

【数14】

$$P_r(a_t = a \mid s_t = s) = \frac{e^{\beta A(s, a)}}{\sum_b e^{\beta A(s, b)}}$$

(式 14)

ただし、 $A(s, a)$  は行動決定のためのパラメータであり、状態  $s$  における行動  $a$  の取りやすさを表している。また、 $\beta$  は行動のランダムさを表すパラメータである。この時、行動生成器のパラメータ更新量  $\Delta A_{\cdot}(s_t, a_{\cdot t})$  と、外乱生成器のパラメータ更新量  $\Delta A_{\cdot}(s_t, a_{\cdot t})$  は(式8-2)の予測誤差  $\delta_t$  を用いて次式でそれぞれ表される。

★ただし、 $\eta^{\cdot}$ 、 $\eta^{\cdot}$  は学習率を表す。また、 $a_{\cdot t}$ 、 $a_{\cdot t}$  はそれぞれ、時刻  $T$  における行動生成器と外乱生成器の行動を表す。

【0026】モデル非依存の学習方式として、行動価値関数を学習することによって、ロバスト強化学習を実現することができる。つまり、次式で表される行動価値関

20 数の予測誤差  $\delta_t$  を用いて、

【数15】

$$\Delta A_{\cdot}(s_t, a_{\cdot t}) = \eta^{\cdot} \delta_t$$

$$\Delta A_{\cdot}(s_t, a_{\cdot t}) = -\eta^{\cdot} \delta_t$$

★

$$\delta_t = q_t + \alpha \max_{a_{\cdot}} \min_{a_{\cdot}} Q(s_{t+1}, a_{\cdot}, a_{\cdot})$$

$$- Q(s_t, a_{\cdot t}, a_{\cdot t}) \quad (\text{式 15})$$

行動価値関数の更新量  $\Delta Q(s_t, a_{\cdot t}, a_{\cdot t})$  は次式のようになる。

$$\Delta Q(s_t, a_{\cdot t}, a_{\cdot t}) = \eta^{\cdot} \delta_t$$

ただし、 $\eta^{\cdot}$  は学習率、 $\alpha$  は評価の減衰率である。

☆【0027】次式で示す確率分布  $P_r$  に従って、状態  $s$  において、外乱生成器の行動  $a_{\cdot t}$  をすべての行動生成器の行動  $a_{\cdot t}$  に関して決定する。

30

☆ 【数16】

$$P_r(a_{\cdot t} = a_{\cdot} \mid a_{\cdot t} = a_{\cdot}, s_t = s)$$

$$= \frac{e^{-\beta Q(s, a_{\cdot}, a_{\cdot})}}{\sum_{b_{\cdot}} e^{-\beta Q(s, a_{\cdot}, b_{\cdot})}} \quad (\text{式 16})$$

この場合、確率分布  $P_r$  に従うことで、小さい行動価値を持つ外乱生成器の行動  $a_{\cdot}$  を高い確率で選択することになる。これによって、目的とする課題の達成にとって外乱生成器が最悪の外乱を生成するようになる。ただし、 $\beta$  は行動のランダムさを表すパラメータである。

40 【0028】次に、次式で表される確率分布  $P_r$  に従って、状態  $s$  において、行動生成器の行動  $a_{\cdot t}$  をすでに決定した外乱生成器の  $a_{\cdot}$  に対する行動  $a_{\cdot t}$  を用いて決定する。

【数17】

$$P_{ru}(a_{\tau+1}=a_{u,i} | s_{\tau}=s) = \frac{e^{\beta Q(s, a_{u,i}, a_{w,\tau+1})}}{\sum_{b_{u,i}} e^{\beta Q(s, b_{u,i}, a_{w,\tau+1})}}$$

(式17)

ただし、行動生成器が行動  $a_{u,i} = a_{u,i}$  を選択した時、外乱生成器は行動  $a_{w,\tau} = a_{w,\tau}$  を選択する。この場合、確率分布  $P_r$  に従うことで、大きい行動価値を持つ行動生成器の行動  $a_{u,i}$  を高い確率で選択することになる。これによって、目的とする課題の達成にとって行動生成器が最高の行動出力を生成するようになる。ただし、 $\beta$  は行動のランダムさを表すパラメータである。

\* {0029} 離散系での行動決定方法:

(モデル依存の場合) 状態  $X_{\tau}$  において、行動生成器の行動が  $u$  で、外乱生成器の行動が  $w$  であり、その結果状態  $X_{\tau+1}$  にたどり着いたとする。そのときに得られる評価信号  $q_{\tau+1}$  とすれば、環境のモデルを用いて、行動生成器の行動  $u_{\tau}$  は

\* {数18}

$$u_{\tau} = \operatorname{argmax}_u \sum_{X_{\tau+1}} P(X_{\tau+1} | X_{\tau}, u, w)$$

$$[q_{\tau+1} + \alpha V(X_{\tau+1})]$$

外乱生成器の行動  $w_{\tau}$  は

$$w_{\tau} = \operatorname{argmin}_w \sum_{X_{\tau+1}} P(X_{\tau+1} | X_{\tau}, u, w)$$

$$[q_{\tau+1} + \alpha V(X_{\tau+1})]$$

となる。

ただし、 $\alpha$  は評価の減衰率、 $P(X_{\tau+1} | X_{\tau}, u, w)$  は状態  $X_{\tau}$  において行動生成器が行動  $u$  を出力し外乱生成器が外乱  $w$  を出力した時、状態  $X_{\tau+1}$  に遷移する確率。

【0030】連続系での行動決定方法:

(モデル非依存の場合) 行動生成器の行動を  $u(t)$ 、外乱生成器の行動を  $w(t)$  とすると、それぞれ、

$$u(t) = A(x(t); v^{Au}) + n_u(t)$$

$$w(t) = A(x(t); v^{Aw}) + n_w(t)$$

\* {数19}

30%のように表される。ただし、 $n_u(t)$ 、 $n_w(t)$  は探索のためのノイズ入力を表す。それぞれの行動は、パラメータ  $v^A = \{v_1^A, v_2^A, \dots, v_n^A, \dots\}$  を持つ近似器  $A(x(t); v^A)$  として実現され、その手段としては、線形モデル、多項式モデルおよび多層神経回路網などを用いることができる。また、それぞれのパラメータは、前述の予測誤差信号  $\delta(t)$  を用いて以下のように更新する。

$$\frac{d v_i^{Au}}{d t} = \eta_u^A \delta(t) n_u(t) \frac{\partial A(x(t); v^{Au})}{\partial v_i^{Au}}$$

$$\frac{d v_i^{Aw}}{d t} = - \eta_w^A \delta(t) n_w(t) \frac{\partial A(x(t); v^{Aw})}{\partial v_i^{Aw}}$$

ただし、 $\eta_u^A, \eta_w^A$  は学習率である。

(モデル依存の場合) 環境のモデルを用いることが可能な場合は、状態評価器の勾配を用いて、モデル非依存性

の場合に比べて効率的に学習を行うことができる。ここ \* [数20]  
で、環境モデルと報酬モデルを次式で表す。 \*

$$\frac{dx}{dt} = f(x) + g_1(x)w + g_2(x)u \quad (\text{式18-1})$$

$$r(x, u) = Q(u) - u^T R(x) u \quad (\text{式18-2})$$

$$s(w) = \gamma^2 w^T w \quad (\text{式18-3})$$

すると、評価関数の勾配と、環境モデルから得られる入 \* 式で表される。

力ゲイン  $g_1(x)$ ,  $g_2(x)$  を用いて、行動生成器の最 10 [数21]

適出力  $u_{**}$  と、外乱生成器の最適出力  $w_{**}$  はそれぞれ次※

$$\text{評価関数の勾配: } \frac{\partial V}{\partial x} \quad (\text{式19-1})$$

$$u_{**} = \frac{1}{2} R(x)^{-1} g_2^T(x) \left( \frac{\partial V}{\partial x} \right)^T \quad (\text{式19-2})$$

$$w_{**} = -\frac{1}{2\gamma} g_1^T(x) \left( \frac{\partial V}{\partial x} \right)^T \quad (\text{式19-3})$$

入力ゲイン  $g_1(x)$ ,  $g_2(x)$  は必ずしも既知ではなくても、状態評価の学習と同時に環境モデルを学習することによって求めることができる。

【0031】具体例1：単振り子の振り上げ

図4の様な単振り子の制御にロバスト強化学習を適用し、学習された制御器を用いて単振り子の振り上げを行った例を示す。単振り子は質量  $m=1$  [kg]、長さ  $L$  [m] で、状態変数は  $x = (\theta, d\theta/dt)$  であり、振り子の角度と角速度で表す。制御指令  $u = \tau$  は振り子の回転軸中心での駆動トルクである。したがって、(式 18-1, 式18-2, 式18-3) との対応を考えると、振り子の運動方程式を構成するそれぞれの関数は以下のように与えられる。

[数22]

$$f(x) = \begin{pmatrix} \dot{\theta} \\ \frac{g}{L} \sin(\theta) - \mu \dot{\theta} \end{pmatrix}$$

$$g_1(x) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$g_2(x) = \begin{pmatrix} 0 \\ \frac{1}{mL^2} \end{pmatrix}$$

$$Q(x) = \cos(\theta) - 1$$

$$R(x) = 0.08$$

$$\gamma = 0.45$$

$\dot{\theta}$  は、文章中では  $d\theta/dt$  で表示する。

40

よって、評価信号は以下の関数で表される。

【0032】

$$q(t) = \cos(\theta) - 1 - 0.08u^2 + \gamma^2 w^2$$

状態変数が  $x = (\theta, d\theta/dt)$  であり、学習時の振り子の質量が  $m=1$  [kg]、長さ  $L=1$  [m] であることから、(式19-1, 式19-2, 式19-3) より、行動生成器および外乱生成器は下記のごとくなる。

[数23]

$$u_{\theta} = 6.25 \frac{\partial V}{\partial \theta}$$

$$w_{\theta} = - \frac{1}{0.405} \frac{\partial V}{\partial \theta}$$

$$\frac{\partial V}{\partial \theta} \text{ は}$$

学習によって獲得される評価関数  $V(x)$  を  $\theta$  で偏微分することで得られる。

ここで、状態変数は連続であるので、状態評価器の関数近似手段として多層神経回路網を用いた。

【0033】このようにして、学習した行動生成器を、非線形ロバスト制御器として採用し、長さ  $L=0.5, 1.0, 2.0$  [m] の3種の異なる長さを持つシステムに適用した。また、従来の強化学習を用いて学習した制御器に対しても同様の実験を行った。その結果を以下に示す。

【0034】図4(b)および図5において、実線で図示するように、全ての環境において、非線形ロバスト制御器は単振り子の振り上げに成功している。ただし、グラフの縦軸は振り子の回転角を、横軸は時間を表してい\*

$$da/dt = \dot{a}$$

$$d\dot{a}/dt = (1/(M+m))$$

$$((F - (M+m)g \sin \theta) \cos \theta - \mu \dot{a})$$

なお、文章中は  $\dot{a}$  は  $da/dt$  で表示する。

ただし、状態変数は  $x = (a, da/dt)$  であり、 $a$  はアクチュエータ11の水平方向の位置を、 $da/dt$  は速度を表す。また、 $F$  はアクチュエータ11が与える力であり、 $M$  は搬送物12の質量、 $m$  はアクチュエータ11の質量である。そして、勾配  $\theta$  は、水平位置が  $a$  の場合には、 $\theta = \arctan(\cos(\pi a))$  となる。

\*る。一点鎖線は振り上がった状態を示している。そして、実線は、一点鎖線で示す直線に収束しているため、振り上げに成功していることが分かる。

【0035】一方、破線で示すように、通常の強化学習で学習した従来の制御器は、学習時に用いた環境と同一の環境以外では振り子を振り上げることができていない。図4(b)に示す様に、学習時の環境(振り子の長さ  $L=1.0$  [m])で振り上げを行うと、振り上げ軌道が一点鎖線に収束していることが分かるが、図5に示したように、学習時の環境以外(振り子の長さ  $L=0.5$  [m],  $2.0$  [m])の環境下で振り上げを行うと、振り上げ軌道は一点鎖線に収束しておらず、振り上げができていないことが分かる。

【0036】具体例2：非線形力場における荷物の搬送  
ここでは、図6の様な直動アクチュエータ11に搬送物12を載せて運搬することを考える。ただし、勾配のために制御対象に非線形性があり、また、アクチュエータ11を小型化するために、大きな出力が出ないような状況を想定する。制御対象の運動方程式は、勾配を  $\theta$ 、摩擦係数  $\mu=0.01$ 、重力加速度  $g=9.8$  [m/s<sup>2</sup>] とすると次式で表される。

【数24】

【0037】したがって、(式18-1、式18-2、式18-3)との対応を考えると、運搬用アクチュエータ11の運動方程式を構成するそれぞれの関数は以下の

【数25】

$$\begin{aligned}
 f(x) &= \begin{pmatrix} \ddot{a} \\ -g \sin \theta \cos \theta - \frac{1}{M+m} \mu \ddot{a} \end{pmatrix} \\
 g_1(x) &= \begin{pmatrix} 0 \\ 1 \end{pmatrix} \\
 g_2(x) &= \begin{pmatrix} 0 \\ \frac{\cos \theta}{M+m} \end{pmatrix} \\
 Q(x) &= \begin{cases} 1.0 & (\text{if } 0.4 \leq x \leq 0.6) \\ -0.5 & (\text{otherwise}) \end{cases} \\
 R(x) &= 0.02 \\
 \gamma &= 0.2
 \end{aligned}$$

ここで、また、(式6)に対応する報酬は以下の関数を\* \*用いた。

$$\begin{aligned}
 q(t) &= 1.0 - 0.02u^2 + \gamma^2 w^2 \quad (\text{if } 0.4 \leq a \leq 0.6) \\
 &= -0.5 - 0.02u^2 + \gamma^2 w^2 \quad (\text{otherwise})
 \end{aligned}$$

【0038】状態変数が $x = (a, da/dt)$ であり、学習時のアクチュエータ11の質量が $m = 1$  [kg]、搬送物12の質量 $M = 1$  [kg]であることから、(式19-2) (式19-3)より、行動生成器および外乱生成器は下記のごとくなる。

【数26】

$$u_{op} = 12.5 \cos \theta \frac{\partial V}{\partial \dot{a}}$$

$$w_{op} = -\frac{1}{0.08} \frac{\partial V}{\partial \dot{a}}$$

$$\frac{\partial V}{\partial \dot{a}} \text{ は}$$

学習によって獲得される評価関数 $V(x)$ を

$\dot{a}$ で偏微分することで得られる。

ここで、状態変数は連続であるので、状態評価器の関数近似手段として多層神経回路網を用いた。

【0039】そして、ある目標地点(図7において一点鎖線で図示する)に移動させることを学習した。なお、図7では、縦軸にアクチュエータ11の位置、横軸に時間を取っている。

【0040】このようにして、学習した行動生成器を、非線形ロバスト制御器として採用し、学習時と同じ質量( $M = 1$  [kg])の搬送物12を載せた場合と、学習時よりも重い搬送物12 ( $M = 3$  [kg])を載せた場合とで、15 [m]離れた地点から目標地点まで搬送するシミュレーション実験を行った結果を図7(a)および図7(b)に実線でそれぞれ示した。図7(a)および

(b)の両方の実線は、アクチュエータ11の軌道が目標地点を示す一点鎖線に収束していることから、搬送物12の質量が、 $M = 1$  [kg]、 $M = 3$  [kg]の両方の場合で搬送を行うことができることが分かる。

【0041】一方、通常の強化学習で学習した従来の制御器に対しても、同様の実験を行い、その結果を図7に破線で図示した。図7(a)の破線で図示したように、搬送物12の質量が $M = 1$  [kg]の場合には、アクチュエータ11の軌道が、目標地点をしめす一点鎖線に収束していることから、目標地点への搬送に成功していることが分かる。しかし、図7(b)の破線で図示したように、搬送物12の質量が $M = 3$  [kg]の場合には、アクチュエータ11の軌道が、目標地点をしめす一点鎖線に収束していないことから、目標地点への搬送ができていないことが分かる。この様に、本発明のロバスト強化学習方式を用いて獲得した制御器は、搬送物12の質量にばらつきがある場合でも、ある程度の範囲内で対応することができる。

【0042】以上、本発明の実施の形態について説明したが、本発明の趣旨の範囲内で種々の形態を実施することが可能である。

【0043】

【発明の効果】以上述べた如く、本発明によれば、目標の達成度に応じた報酬に、外乱に耐えうることに応じた報酬を加味した報酬信号である評価信号を生成し、現在の状態から将来に向けて得られる評価信号の和の期待値を最大化するべく行動生成器は学習し、一方、外乱生成器は前記評価信号の和の期待値を最小化すべく学習するので、環境モデルが未知の場合にも対応することができる。さら

に、非線形の制御対象あるいは環境にも適用可能である。

【図面の簡単な説明】

【図1】本発明の実施の形態のロバスト強化学習方式に用いる回路のブロック図である。

【図2】従来の強化学習方式に用いる回路のブロック図である。

【図3】 $H_\infty$ 制御理論を説明するための制御対象と制御器との制御系のブロック図である。

【図4】第1具体例の説明図で、(a)が概略図、(b)が1mの長さの振り子の角度変化のグラフである。

【図5】第1具体例の角度変化のグラフで、(a)が0.5mの長さの振り子を制御対象として用いた場合のグラフ、(b)が2mの長さの振り子を制御対象として用いた場合のグラフである。

\*【図6】第2具体例の概略図である。

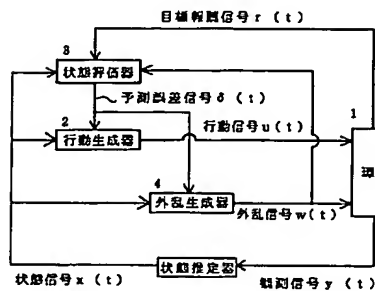
【図7】第2具体例の直動アクチュエータの位置変化のグラフで、(a)が1kgの質量の搬送物を用いた場合のグラフ、(b)が3kgの質量の搬送物を用いた場合のグラフである。

【符号の説明】

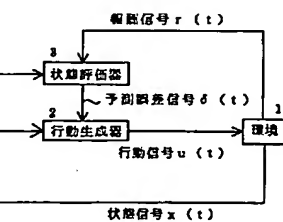
- q(t) 評価信号
- r(t) 目標報酬信号
- s(t) 外乱報酬信号
- u(t) 行動信号
- w(t) 外乱信号
- 1 環境
- 2 行動生成器
- 3 状態評価器
- 4 外乱生成器

\*

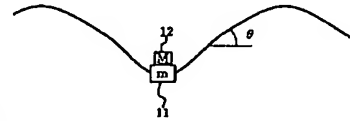
【図1】



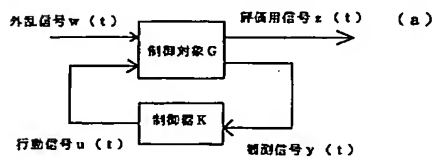
【図2】



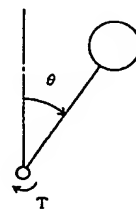
【図6】



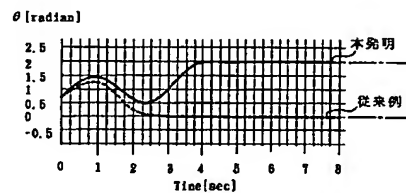
【図3】



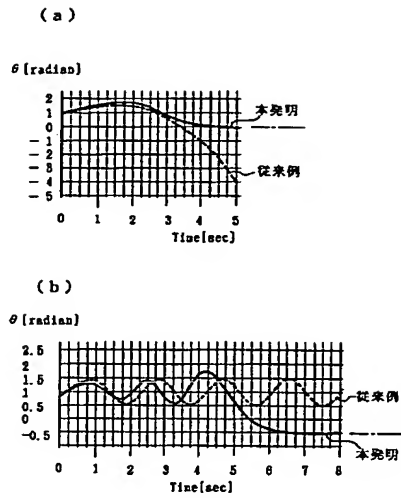
【図4】



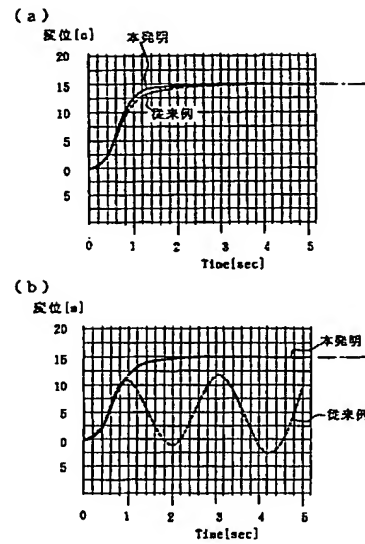
(b)



【図5】



【図7】



フロントページの続き

(72)発明者 銅谷 賢治  
京都府相楽郡精華町光台 7-2-1-5  
- 201

Fターム(参考) 5H004 GA07 GA15 GA17 JA13 JB22  
KC09 KC18 KC28 KD42 KD62